

(51) Int.Cl. ⁶	識別記号	庁内整理番号	F I	技術表示箇所
G 0 6 F 17/30				
12/00	5 4 5 A	8944-5B		
// G 0 6 F 9/44	5 5 0 N	9193-5B		
		9194-5L		
			G 0 6 F 15/ 401	3 1 0 C

審査請求 未請求 請求項の数 6 O L (全 12 頁)

(21) 出願番号 特願平5-327352

(22) 出願日 平成5年(1993)12月24日

(71) 出願人 000005108

株式会社日立製作所

東京都千代田区神田駿河台四丁目6番地

(72) 発明者 前田 章

神奈川県川崎市麻生区王禅寺1099番地 株

式会社日立製作所システム開発研究所内

(72) 発明者 芦田 仁史

神奈川県川崎市麻生区王禅寺1099番地 株

式会社日立製作所システム開発研究所内

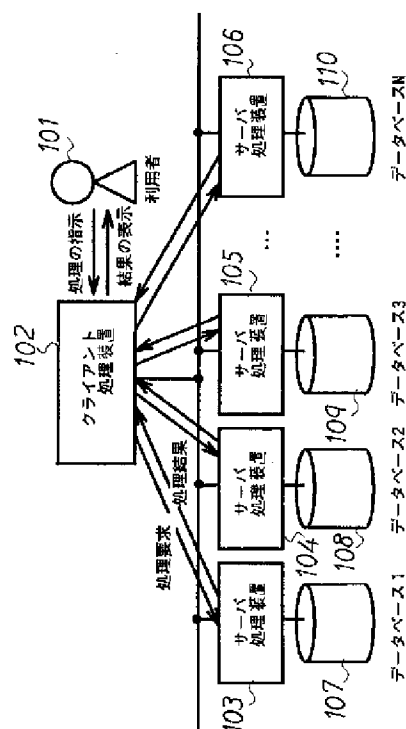
(74) 代理人 弁理士 磯村 雅俊

(54) 【発明の名称】 データ加工システム

(57) 【要約】

【目的】 複数のデータベースに分割して格納されたデータに対して加工処理を行う際に、複数のデータベースに設けられた処理装置を並列に活用することにより高速な処理を実現し、かつ複数のデータベースを接続するネットワークの負荷を軽減し、かつ利用者にとってデータの物理的な配置を意識せずにデータ加工処理を実行することを可能にすること。

【構成】 クライアント処理装置102から複数のサーバ処理装置103～106に分類処理要求を送出し、それを受けた複数のサーバ処理装置は、分類処理を並列に実行し、処理結果をクライアント処理装置に伝達する。処理結果を受けたクライアント処理装置は、それらを合成し、合成した結果を用いてルールインダクションなどの処理を実行する。



【特許請求の範囲】

【請求項 1】 クライアント処理装置と、該クライアント処理装置に接続され、それぞれがデータベースに接続された複数のサーバ処理装置からなるデータ加工処理システムにおいて、

処理対象となる 1 つ以上のテーブル形式のデータを上記データベースの各々に格納する格納手段と、

上記テーブル形式のデータに含まれるデータ項目の値を分類するための分類規則を指定する分類規則指定手段と、

上記分類規則指定手段により指定された分類規則を上記サーバ処理装置に伝達する第 1 の伝達手段と、

上記サーバ処理装置において、上記伝達された分類規則にしたがってそれぞれのデータベースに格納されているデータに対して分類処理を実行する分類実行手段と、

上記分類実行手段によって分類された結果を上記クライアント処理装置に伝達する第 2 の伝達手段と、

上記クライアント処理装置に伝達された分類結果を合成して上記処理対象であるテーブル形式のデータに対する分類結果を生成する手段と、

上記生成された分類結果を用いて、データ項目間の関係を分析する分析手段と、

上記分析手段によって分析された結果を出力する出力手段とからなることを特徴とするデータ加工システム。

【請求項 2】 請求項 1 記載のデータ加工システムにおいて、上記複数のサーバ処理装置は同時並列に処理を行うことを特徴とするデータ加工システム。

【請求項 3】 請求項 1 または 2 記載のデータ加工システムにおいて、上記分類規則は、実数または整数の値を持つデータ項目に対しては上限値および下限値を指定し、文字列としての値を持つデータ項目に対しては文字列の変換を指定するものであることを特徴とするデータ加工システム。

【請求項 4】 請求項 1 ～ 3 記載のデータ加工システムにおいて、上記分類規則には、テーブル形式のデータの各レコードを一意に識別する項目が含まれることを特徴とするデータ加工システム。

【請求項 5】 請求項 1 ～ 4 記載のデータ加工システムにおいて、上記分析されたデータ項目間の関係を、1 つまたは複数のルールとして出力することを特徴とするデータ加工システム。

【請求項 6】 請求項 1 ～ 4 記載のデータ加工システムにおいて、上記分析されたデータ項目間の関係を、分類木として出力することを特徴とするデータ加工システム。

【発明の詳細な説明】

【0001】

【産業上の利用分野】 本発明は、物理的に複数のデータベースに分割して格納された数値または記号で表現された情報の集まりを加工し、利用者にとって有用な表現に

変換して出力するデータ加工システムに関し、特に、高速な処理を実現し、かつ利用者にデータが物理的に分割されて格納されていることを意識させずに処理を実行することができるデータ加工システムに関する。

【0002】

【従来の技術】 最近の情報処理技術にはめざましいものがあり、電子計算機内に蓄積されるデータ量や処理すべきデータ量は年々増大しており、特にネットワーク化が進むにつれて、オンラインシステムを中心にこの傾向はますます顕著になってきている。現在では、そのデータ量はギガバイト（＝10⁹乗）、レコード数にして100万件を超えるものも珍しくない。一般に、電子計算機内に蓄積されたデータそのものは単なる数値や記号の集合に過ぎず、そのままでは利用することができない。そこで、データの有効活用を図るために、このデータの集合を有用な情報に変換して利用者に提供するための技術がいろいろ提案されている。その中の代表的な技術として、

（1）回帰分析や重相関分析といった統計的な手法（以下、従来技術 1 という）

（2）ニューロモデルやファジィモデルを用いて、データ間の相互関係を学習させる方法（以下、従来技術 2 という）

（3）ルールインダクションなど、知識獲得手法を用いる方法（以下、従来技術 3 という）

などが従来から公知である。

【0003】 上記従来技術 1 に関しては、広く一般的に知られている周知の手法であり、ここで詳しく述べることはしない。上記従来技術 2 に関して、ニューロモデルの学習に関しては多くの文献に記載されているが、これは基本的には多入力多出力の非線形関係を学習によりモデル化するものである。また、非線形関係の学習という意味では、入出力関係をファジィモデルで表現し、ニューロと同様なアルゴリズムを用いて学習によりメンバシップ関数の形状を調節する方法が知られている。ニューロまたはファジィいずれの方法にしても、入出力関係のあるパラメータを含むモデルで表現しておき、そのパラメータを学習により決めることによりモデルを構築するものである。

【0004】 上記従来技術 3 に関しては、例えば、Springer-Verlag の「Machine Learning」, 463～482 ページに記載された“Learning Efficient Classification Procedures and Their Application to Chess End Games”という J. Ross Quinlan の論文がある。この論文には、データから分類のための決定木を自動的に作成する ID3 というアルゴリズムが述べられている。また、ルールインダクション手法として、日立クリエイティブワークステーション 2050 マニュアル、「ES/TOOL/W-R I 解説/操作」の第 23 頁～第 53 頁にデータ間に存在する関係をルールの形で表現

する手法が詳細に説明されている。

【0005】次に、ID3による決定木作成例を説明しておく。まず、全体の事例の集合Cは、

C = {short, blond, blue, : +, short, dark, blue : -, tall, dark, brown : -, tall, blond, brown : -, tall, dark, blue : -, short, blond, brown : -, tall, red, blue : +, tall, blond, blue : +}

で与えられているものとする。各事例は身長（値はshort, tall）、髪の色（値はblond, dark, red）、目の色（blue, brown）という3つの属性を持ち、さらに+と-という2つのクラスの何れかに属している。ID3は、どの属性をどういう順番で判定すれば、事例を+または-というクラスに分類できるかという問題を、決定木を生成することにより解決しようとする。決定木の生成は、どの属性から判定すれば最も判定後の情報量が最小になるか（すなわち、どの属性による判定が最も多くの情報量を持つか）という基準で属性を選んでいく。図8に、上の事例の場合に得られる決定木の例を示す。この決定木の意味するところは次のようなものである。まず、

（a）事例を髪の色で分類せよ。

（a1）髪の色がdarkならば-である。

（a2）髪の色がredならば+である。

（a3）髪の色がblondならば、次に目の色で分類せよ。

（a31）目の色がblueならば+である。

（a32）目の色がbrownならば-である。

これまで、いろいろなID3の変形アルゴリズムが提案されているが、基本的にはある基準にしたがって判定の順序を決定するものである。

【0006】上記従来技術3は、上記従来技術1および2と異なり、データを加工した結果がルールインダクションの場合にはルールとして、またID3の場合には決定木として、陽に表現された形で得られることから、蓄積されたデータに含まれている関係を利用者が発見するのに特に有効である。これらの技術は大量データの有効利用技術として用いられることが多くなっている。

【0007】上記のようなデータの利用技術と並行して、大量のデータを高速に検索することを主目的とした並列データベース技術の開発が進んでいる。このような並列データベース技術は、例えば、「日経エレクトロニクス1993年7月19日号」の第91頁～第106頁に記載された“並列マシン向けDBMS技術－90年代半ばの実用化めざす”という論文に解説されている（以下、従来技術4という）。また、利用者のインタフェースとして表計算ソフトを用いてデータベースの検索を行うソフトウェア技術が、例えば、「日経コンピュータ1993年7月12日号」の第65頁～第75頁に記載された“サーバのRDBを直結－基幹データを表計算ソフトへ”と題する論文に解説されている（以下、従来技術5という）。この論文に記載されているものでは、サー

バ処理装置と、利用者が直接使用するクライアント処理装置をネットワークで接続し、クライアント処理装置では表計算ソフトが提供するインタフェースを用い、データベースの検索要求が発生した場合にはSQL（Structured Query Language）などで記述された検索命令をサーバ処理装置に送り、サーバ処理装置がその要求にしたがって検索を実行し、検索結果をクライアント処理装置に転送している。クライアント処理装置はサーバ処理装置から転送された検索結果を表計算の形に整形して利用者に対して表示する。このような仕組みにより、少なくとも検索や簡単な統計量の算出（平均値や分散など）に関しては、利用者は全く分散データベースを意識せずに、全てデータが手元（クライアント処理装置）にあるかのように作業することができる。

【0008】

【発明が解決しようとする課題】上記従来技術1および2は、電子計算機内に蓄積されたデータを回帰分析、統計処理または学習処理によって加工する技術であるが、加工した結果を利用者がどのように利用するかなど、その利用方法については特に考慮されていなかった。それに対して、従来技術3のルールインダクションやID3アルゴリズムなどの知識獲得手法を用いる方法では、加工された結果は利用者にとって有用なものになりうるが、この技術は主にどのようなデータ加工手段を用いれば有用な情報を得ることができるかという点に関するものであり、従来技術4で扱っているような並列データベースの上でどのように実行するかという点については、これまで検討されることがなかった。

【0009】一方、従来技術4の並列データベース技術は、いわゆるオンラインシステムにおけるトランザクション処理と、オフライン処理におけるデータ検索処理を、いかに並列ハードウェアを用いて高速に実行するかに主眼がおかれたものであり、上記従来技術1～3のように大量のデータを高度に加工するような処理を、ハードウェアの並列性を生かして高速に実行するような手段についてはこれまで知られていない。したがって、従来技術3の知識獲得手法を従来技術4の並列データベースに適用することによる生じる効果は、単に、データの検索の高速化という並列データベース本来の効果に限られていた。このように、従来は、ハードウェアが本来持っている並列能力を生かした方法については知られておらず、したがって、従来技術3の知識獲得手法を十分高速に実行することができないという問題点があった。本発明の第1の目的は、上記の問題点を解決し、並列・分散データベース上で高速に実行可能なデータ加工方法および装置を提供することにある。

【0010】また、一般に、新規なデータを抽出収集して蓄積するには多大な手間、時間、およびコストがかかるため、既存のデータ資産をいかに有効に利用するかが重要な技術的課題となっている。例えば、同一の会社

の中でも、異なる部門毎にそれぞれ別個の情報をまったく別のデータベースとして構築していることがある。この場合、データベースは物理的にも遠く離れ、まったく別種のハードウェア上に、まったく別種のソフトウェアを用いて構築されていることさえある。上記の従来技術3においては、データは処理を実行する処理装置上に存在すると仮定していた。すなわち、従来技術3を適用するためには、分散したハードウェア上にある関連データを切り出して、1つの処理装置上に集めてくる処理が必要になる。現在では、処理装置はネットワークで接続されていることが多いが、それでも関係するデータを集めたデータベースをデータ加工を試みる毎に定義して構成するというのは、そのために時間を要し、また、ネットワークの負荷としても、また人的作業の点からも、非常にコストがかかるという問題点があった。上記従来技術3を利用する側から見た場合、利用者は様々な目的を達成するためにデータ加工を行うのであって、そのためのデータの準備に多くの時間、コストをかけるというのは非効率的であり、さらに、これらの技術の有効性を十分に生かしきれていないという問題点があった。

【0011】上記問題点を、データベース検索の場合に解決しようとする技術が上記従来技術5である。しかしながら、この従来技術5の対象とするのは、あくまでも通常の表計算で用いるような処理だけであり、上記従来技術3で対象としているような高度なデータ加工処理に対しては対策がなされていなかった。例えば、ルールインダクションによるルール抽出処理を実行しようとする、クライアント処理装置上の処理でデータを参照するたびにサーバ処理装置へのデータ要求が発生するため、ネットワークの負荷は非常に重くなり、ほとんど実用的な処理速度にはならないという問題点がある。特に、データベースが遠く離れた場所にある場合、例えば、あるデータベースが東京に、また別のデータベースが大阪に存在するような場合には、この問題点はより顕著になって現れる。本発明の第2の目的は、上記の問題点を解決し、利用者にとっては物理的なデータの存在場所を意識することなく、ネットワークで接続されたサーバ処理装置上のデータを用いて、ルールインダクションなどの手法を用いてデータ間の関係を利用者にとって有用な形に加工することができ、またこれらの処理をネットワークに接続された処理装置の処理能力を有効に活用することにより、高速に実行する方式および装置を提供することにある。

【0012】

【課題を解決するための手段】本発明は、上記の問題点を解決するために、クライアント処理装置と、該クライアント処理装置に接続され、それぞれがデータベースに接続された複数のサーバ処理装置からなるデータ加工処理システムにおいて、処理対象となる1つ以上のテーブル形式のデータを上記データベースの各々に格納する格

納手段と、上記テーブル形式のデータに含まれるデータ項目の値を分類するための分類規則を指定する分類規則指定手段と、上記分類規則指定手段により指定された分類規則を上記サーバ処理装置に伝達する第1の伝達手段と、上記サーバ処理装置において上記伝達された分類規則にしたがってそれぞれのデータベースに格納されているデータに対して分類処理を実行する分類実行手段と、上記分類実行手段によって分類された結果を上記クライアント処理装置に伝達する第2の伝達手段と、上記クライアント処理装置に伝達された分類結果を合成して上記処理対象であるテーブル形式のデータに対する分類結果を生成する手段と、上記生成された分類結果を用いて、データ項目間の関係を分析する分析手段と、上記分析手段によって分析された結果を出力する出力手段とを有している。

【0013】

【作用】本発明は、上記各手段、特に、テーブル形式のデータに含まれるデータ項目の値を分類するための分類規則を指定する分類規則指定手段と、分類規則指定手段により指定された分類規則を上記サーバ処理装置に伝達する第1の伝達手段と、サーバ処理装置において上記伝達された分類規則にしたがってそれぞれのデータベースに格納されているデータに対して分類処理を実行する分類実行手段と、分類実行手段によって分類された結果をクライアント処理装置に伝達する第2の伝達手段と、クライアント処理装置に伝達された分類結果を合成して処理対象であるテーブル形式のデータに対する分類結果を生成する手段と、生成された分類結果を用いてデータ項目間の関係を分析する分析手段とを有することによって、時間のかかる処理を並列に実行することができ、かつ利用者には並列処理を全く意識させることなく、また大量のデータ自身を転送することがなく、目的の処理を実行することができる。また、ルールインダクションに適用した場合、ルールインダクションの処理自身はクライアント処理装置上で実行されるが、クライアント処理装置上には各データベース上で分類された結果だけを保持していればいいので、ネットワークの負荷を増大させることなく、目的の処理を実行することができる。さらに、利用者は使用しようとするデータが物理的にどのデータベース上に存在するかということは全く意識することなく処理を実行し、かつ適切な処理結果を得ることができる。

【0014】

【実施例】まず、本発明の実施例の概略を説明する。図1は本発明を実施するための全体構成である。図1において、101は本発明のデータベース加工装置を利用する利用者、102はクライアント処理装置、103～106はサーバ処理装置、107～110はそれぞれサーバ処理装置103～106に接続されているデータベース、111はクライアント処理装置とサーバ処理装置を

接続するバスである。複数のデータベース107～110は互いに独立なDBMS（Date Base Management System）で管理されている場合であっても、全くネットワーク上で統一した管理を行う並列DBMSで管理されている場合であってもよい。ここでは前者の場合を例にとって説明する。また、説明の便宜上、データは図2に示したように複数の属性からなる複数のレコード（行に対応）によって記述されており、それら複数のレコードは複数のデータベースに分割されて格納されているものとする（図2では、それぞれK個の属性を有するM個のレコードが3個のデータベースに分割されて格納されることを示している）。

【0015】例えば、クライアント処理装置102では、複数のデータベース107～110に分散されて格納されているデータに対して、仮想的に一つのまとめ上げられた表のイメージで表示することができる（従来技術5）。このとき、利用者は、その表イメージのデータに対して、例えば、ルールインダクションによる分類ルールの抽出を行うものとする。複数のデータベース上に蓄積されているデータは一般に膨大なものであるから、クライアント処理装置に全てのデータを転送することは処理時間とコストの点から一般には行うことができない。

【0016】そこで、本発明では、クライアント処理装置上で動作するルールインダクション処理が、ある属性または属性値による分類結果を求める処理要求を複数のサーバ処理装置に送信する。処理要求の形式は、例えば、一般によく知られているSQL（Structured Query Language）文の形式でもよい。複数のサーバ処理装置は、クライアント処理装置から要求された処理をそれぞれ並列に実行し、処理結果をクライアント処理装置に伝達する。ルールインダクション処理の場合、この処理結果は、クライアント処理装置によって指定された属性または属性値によって各データベース中のレコードを分類した結果を、例えば、レコード番号と分類結果の対応表の形で表現されている。クライアント処理装置は、複数のサーバ処理装置から伝達された処理結果をまとめ上げることにより、以後のルールインダクションに必要な情報、すなわちそれぞれの属性または属性値による全レコードの分類表をクライアント処理装置上に作成することができる。したがって、データベースに蓄積された情報自身をクライアント処理装置に転送することなしに、分類結果だけをクライアント処理装置上にもつことによって、ネットワークの負荷を軽減し、高速に処理をすることが可能になる。

【0017】また、このような処理を行わずに、クライアント処理装置が直接ルールインダクションを実行した場合には、ルールインダクション処理の進行に従って複数のサーバ処理装置への検索要求が逐次伝達されることになり、やはりネットワークの負荷が非常に大きくな

り、結果としてルールインダクション処理に要する時間が多大なものになってしまう。本発明で設けた手段により、このようなオーバーヘッドを低減できるので、高速な処理を実現することができる。

【0018】以下、本発明の実施例の詳細な動作を説明する。実施例の説明にあたって、まず、ルールインダクションの動作原理について説明しておく。詳細については、例えば、上述した文献、日立クリエイティブワークステーション2050マニュアル、「ES/TOOL/WRI解説/操作」第23頁～第53頁に解説されている。図3はルールインダクションの入力となる事例テーブルの例である。ルールインダクションの入力となる事例テーブル200において、各行がそれぞれ一つの事例に対応している。この例では、クレジットの申し込みがそれぞれ一件の事例に対応する。それぞれの事例は1つまたは複数の属性をもっている。図3の例では、「銀行の信用」201、「預金残高」202、「現在の負債」203、「クレジット査定」204という4つの属性がある。それぞれの属性は、あらかじめ定められた形式の属性値をもつ。図3の例では、「銀行の信用」201という属性は「有」と「無」という2つの属性値のいずれかの値をもつ。このような属性を「記号属性」と呼ぶ。それに対して「預金残高」202と「現在の負債」203という属性は数値で表される属性、すなわち、数値属性をもっている。一般の事例テーブルでは、記号属性と数値属性が混在する。また、図3の例において、事例の中にはその属性の一部が不明なものも存在する（参照符号205参照）。これを「欠損値」と呼ぶ。

【0019】ルールインダクションの実行にあたっては、まずどの属性を結論属性とするかを指定する。図3の例では、「クレジット査定」を自動的に判別するシステムを作ることを目的として、「クレジット査定」という属性が結論属性として指定されるものとする。一般のルールインダクション処理では結論属性は記号属性でなければならない。さらに、各事例における結論属性の属性値を説明するために入力として扱われる説明属性を指定する。図3の例では、結論属性以外の属性をすべて説明属性として扱うものとする。もちろん、事例テーブル中の一部分だけを説明属性として指定することもできる。

【0020】図3のような事例テーブルにおいて、事例の数（行の数）が大きくなるにつれて、数値属性をそのまま扱うことが困難になってくる。そこで、数値属性を記号属性に置き換える「ラベル付与処理」をルールインダクション処理に先立って行うことがある。ラベル付与処理とは、数値属性をもつ属性に対して、例えば、適当な区間分割を行って、各区間にラベルと呼ばれる記号を割り当てることによって数値を記号に置き換える処理のことを指す。図4はこのようなラベル付与処理を行ってすべての属性を記号属性に置き換えた事例テーブルの例

である。これは、「預金残高」208が1000万円以上を「大」、1000万未満で正のものを「中」、預金残高マイナスのものを「小」とした例である。「現在の負債」に関しては、8000万円以上を「大」、8000万円未満かつ4000万円以上のものを「中」、4000万円未満のものを「小」としてラベル付与した場合である。

【0021】このように、実数または整数の値を持つデータ項目に対してはそれぞれ上限値および下限値を指定することによっていくつかのグループに分類することができる。また、上限値および下限値を指定する際に、各グループが同じ数、ほぼ同じ数になるように分割するようにしてもよい。ラベル付与処理によって、数値属性の細かな差を無視し、大きく分類した記号属性値にしたがって属性間の関係を解析することができる。特に事例数が大きな事例テーブルに対しては、処理時間の点でも出力される分類規則の有用性という点でも有効である。また、文字列としての値を持つデータ項目に対しては、予めいくつかの文字列を決めておき、そのいずれかの文字列に変換することによって適当な数のグループに分類することができる。

【0022】次に、ルールインダクション処理では、結論属性として指定された属性の属性値を、できるだけうまく説明するような説明属性値の組を見つけ、それをルールとして抽出する処理を行う。この処理にはさまざまなオプションがあり、それらのオプションを適当に設定することにより、利用者にとってより有用なルールを抽出することができる。ここではこれらのオプションについては説明しない。詳細は上記文献、日立クリエイティブワークステーション2050マニュアル、「ES/T OOL/W-R I 解説/操作」の第23頁～第53頁に記載されている。

【0023】図5(a)、(b)は、ルールインダクション処理によって抽出されたルールの一例である。ルールインダクションの内部処理についてはさまざまな方式が提案されているが、基本的には、1つまたは複数の属性が特定の属性値（または属性値の組）をもつ事例の数をベースにして、結論属性を最もよく分類する属性値の組を見つける処理からなる。したがって、ルールインダクション処理を実行するに当たっては、各事例の各属性がどの記号属性値をもつかという対応表があれば十分であることが分かる。すなわち、ルールインダクションの初期処理において、このような対応表を準備しておけば、それ以後の処理は全てこの対応表だけに基づいて実行することができる。

【0024】図6にこのような対応表211の例を示す。図6の対応表の各行は、図3の事例テーブルと同様にそれぞれが一つの事例に対応する。図6の各列は、それぞれ属性「銀行の信用」213、「預金残高」214、「現在の負債」215、「クレジット査定」216

に対応している。対応表211の各成分は、対応する事例の対応する属性値がどの値をもつかを表すラベルコードを示している。図6では簡単のため、事例を区別するためにユニークな事例番号（レコード番号）212を付与してある。また、対応表211の各成分の値がどの記号属性値に対応するかは、図7に示したような別の対応表で管理することができる。図7において、(a)は銀行の信用の各成分（無/有）と記号属性値（0/1）の対応表、(b)は預金残高の各成分（小/中/大/不明）と記号属性値（0/1/2/3）の対応表、(c)は現在の負債の各成分（小/中/大）と記号属性値（0/1/2）の対応表、(d)はクレジット査定の各成分（不可/可）と記号属性値（0/1）の対応表である。

【0025】ここで重要なのは、各属性に対する属性値が何種類の値をとるかに応じて、ラベルコードのとり値は一般に数ビットで表現できることである。図3の事例テーブル自身は、例えば、「クレジット査定」の属性に対して、“不可”という文字を表すために4バイト（全角漢字2つ分）のデータ量を割り当てる必要がある。それに対して、図6の対応表のラベルコードは“1”と“0”を識別できるだけでよいから、基本的には1ビットでよい。したがって、ルールインダクションを非常に大量の事例に対して実行する場合、図3または図4の事例テーブルから図6の対応表にあらかじめ変換しておくことにより、効率のよい処理が可能になる。以上、ルールインダクション処理の概要について説明した。ID3などの分類の決定木を作成するアルゴリズムも、上記の対応表の参照だけで実行できることはいうまでもない。

【0026】さて以上の処理は、本発明の手段を用いることによって、図1に示したような環境で実行することができる。すなわち、図1の構成では、ルールインダクション処理を実行するクライアント処理装置と、データの存在するサーバ処理装置が物理的に異なる処理装置であることと、またデータ自身も複数のサーバ処理装置に分散して格納されているから、上記の処理をそのまま単純に実行することはできない。全ての関係するデータをクライアント処理装置上に転送し、クライアント処理装置上で処理を集中的に行うことができれば話は簡単であるが、特に大量のデータを扱う場合や、データベース自身が物理的に遠く離れた場所にあつてデータ転送に大きなコストがかかる場合には実際的ではない。

【0027】以下、本発明の第1実施例を詳細に説明する。第1実施例は、図1および図2の構成においてルールインダクション処理を高速かつ効率的に実行するものである。図8に事例テーブル220の構成例を示す。この例では、事例を区別する「事例番号」221、各事例の属性値として、「氏名」222、「氏名コード」223、「銀行の信用」224、「預金残高」225、「年収」226、「現在の負債」227、「年齢」228、「性別」229、「住所」230、「配偶者の有無」2

31、「扶養家族人数」232、「職業」233、「最終学歴」234、「住居の種類」235、「クレジットカード保有枚数」236、「利用実績」237、「クレジット査定」238からなっている。

【0028】図9に、図8に示した事例テーブル220を複数のデータベースへの分割して格納する方法を示す。本実施例では、図9のように事例番号1～1000はデータベース1に、事例番号1001～2500はデータベース2に、事例番号2501～4500はデータベース3に、というように事例テーブルを行（レコード）で分割して格納するものとする。各サーバ処理装置は、ローカルにどの範囲のレコードをもっているかを管理しており、クライアント処理装置からの検索処理要求を受けると、ローカルなデータの範囲でだけ検索を実行して結果をクライアント処理装置に伝達する機構を設けている。この方法は、並列DBMSでは「shared nothing」方式と呼ばれているものである。もちろん、並列DBMSのように全体を統一的に管理するDBMSが存在せず、単に独立したローカルなDBMSでそれぞれのデータベースが管理されている場合でも、以下の処理は同様である。

【0029】全体処理の流れを図10に示す。まず、ステップ250において、クライアント処理装置（図1参照）は、図8に示された形式で事例テーブルの一部を表示装置に表示する。利用者は表示装置上で、結論属性と説明属性を指定する（ステップ251）。ここで、一般的に考えて明らかに結論属性と関係がない属性は省しておく。この例では、「氏名」222、「氏名コード」223、「住所」230などの属性は結論属性「クレジット査定」と関係がないと考えられるから予め省しておく。ここではそれ以外の全ての属性を説明属性として指定した場合を考える。

【0030】次のステップ252において、クライアント処理装置は、説明属性および結論属性に対して、ラベル付与処理に必要な情報である「ラベルコード対応表」を作成する。図11にこのラベルコード対応表261の例を示す。ラベルコード対応表261において、参照符号262は、「残高」という属性に対して定義された「小」というラベルは属性値が0未満の値に対応し、「中」というラベル値は属性値が0以上1000未満に対応し、「大」というラベル値は属性値が1000以上に对应していることを意味している。ラベルコード対応表261中で値が指定されていない部分は制約が指定されていないものとする。また、欠損値に対応する「不明」というラベル値に対しても、特別のラベルコード値を割り当てるものとする。例えば、ラベルコード対応表261中に参照符号263で示したように、残高不明に対してラベルコード“3”を対応させている。また、参照符号264、265に示したように、記号属性に関しても、属性値とラベルコードの対応関係を指定する。

【0031】次に利用者によってルールインダクション処理の開始が指示される（ステップ253）と、ステップ254において、クライアント処理装置は、複数のサーバ処理装置に対し、図11で示した情報からなる検索要求を送信する。

【0032】ステップ255において、各サーバ処理装置（図1参照）は図11の検索要求を受信すると、それぞれのサーバ処理装置において検索処理を開始する。まずそれぞれの事例データを、図11のラベルコード対応表に基づいて、各事例を図6と同様の事例－属性値対応表に変換する（ステップ256）。図12にその一例を示す。この検索処理では、図11のラベルコード対応表に現れる全ての属性に対して、その属性が記号属性であれば属性値を指定されたラベルコードに置き換え、数値属性であればラベルコード対応表に指定された数値範囲で分類してラベルコードに置き換える処理を行う。各サーバ処理装置は、ローカルに管理している事例テーブルを図12の事例－属性値対応表に変換した後、ステップ257において、その結果を検索結果としてクライアント処理装置に送信する。

【0033】クライアント処理装置は、複数のサーバ処理装置から処理結果を受信し（ステップ258）、それを合成して、事例テーブル全体の事例－属性値対応表をクライアント処理装置上に作成する（ステップ259）。この処理は、図9に示した事例テーブルの分割と全く逆に、それぞれのサーバ処理装置からの検索結果を行方向に結合すればよい。この事例－属性値対応表が作成された後は、ルールインダクション処理はこの表だけを参照してルールを抽出する（ステップ260）。以上の説明は、ルールインダクションに限らず、前述したID3のような属性に基づいて分類の決定木を自動生成する処理にも全く同様に適用することができる。

【0034】以上、本発明の第1の実施例によれば、クライアント処理装置におけるルールインダクション処理の実行に当たって、あらかじめ複数のサーバ処理装置で実行されるDBMSの機能を拡張しておき、クライアント処理装置で必要な検索結果をクライアント処理装置からの要求に従って並列に実行させることが可能になる。複数のサーバ処理装置への検索要求とその結果の取得は、ルールインダクション処理本体の実行に先立って一度だけ実行すればよく、ネットワーク上で流れるデータ量を大幅に削減し、ネットワーク負荷を減少させることが可能になる。また、検索処理自身も並列に実行されるため、全体の処理時間を短縮できるという効果もある。さらに利用者から見れば、あくまでクライアント処理装置上の表イメージのデータに対して直接処理を実行し、結果を得るという作業が可能になるため、事例テーブルが複数のデータベースに分割して格納されていることは全く意識することがなく、より自然に作業を進めることができるという効果もある。

【0035】本発明の第1の実施例では、図2に示すように、事例テーブルを行単位に分割して複数のデータベースに格納することを前提としていた。すなわち、複数のデータベースは、あらかじめ組となって一つの事例テーブルを格納するために使用されていたことができる。ところが実際には、複数の部署が全く独立にデータベースを構築し、それらのデータベースを統合して利用する場合がある。その場合には、事例テーブルの分割は、行単位ではなく、列単位に分割されていると考えることができる。

【0036】本発明の第2の実施例は、列単位で分割された複数のデータベース上のデータを用いてルールインダクションなどの処理を行うものである。図13に、本実施例におけるデータベースの構成を示す。データベースの内容は第1の実施例で説明した図8の構成と同様であるが、家族情報と金融関連情報が別個のデータベースに格納されている。図13(a)は家族情報に関するデータベース300であり、「氏名」301、「氏名コード」302、「年齢」303、「性別」304、「住所」305、「配偶者の有無」306、「扶養家族人数」307、「職業」308、「年収」309、「最終学歴」310、「住居の種類」311という属性からなっている。また、図13(b)は金融関連情報データベース320であり、「氏名」321、「氏名コード」322、「銀行の信用」323、「預金残高」324、「年収」325、「職業」326、「クレジットカード保有枚数」327、「利用実績」328、「クレジット査定」329という属性からなっている。

【0037】これらのデータベースは、当初は全く異なる目的のために構成されたものであってもよく、またハードウェア的にも全く異なる装置上に存在していてもよい。さらには、物理的に遠くはなれた場所に存在していてもかまわない。ただし、これらのデータベースはネットワークに接続されていて、利用者が直接使用するクライアント処理装置からオンラインで利用できなければならない。ここで、2つのデータベースに「氏名」、「氏名コード」、「職業」、「年収」という属性が共通して存在する。このうち「氏名コード」をレコードを特定するためのキー属性と考える。2つのデータベースの整合性を管理していない場合は、「職業」や「年収」という属性の値は、2つのデータベースで異なることもありうる。そこで、利用者は図8のテーブルの形でデータを扱う際に、各項目が物理的にどのデータベースのどういう項目に関係づけられているかを指定しなければならない。

【0038】図14にこの関係を定義する属性関係づけテーブル330の例を示す。図14のテーブル330は、属性、データベース種別、項目名から構成され、説明属性・結論属性・キー属性などの各属性がそれぞれのデータベース中の(データベース種別)どの属性(項

目名)に対応しているかを表している。例えば、図14のテーブル330の例では、「銀行の信用」、「残高」という属性は、金融関係データベースの項目名「信用」、「預金残高」に対応させ、「氏名」、「年収」という属性は家族情報データベースの項目名「氏名」、「年収」に対応したものであることを示している。キー属性はレコードを特定するものであるから、当然使用する全てのデータベースに属性として含まれている必要がある。このテーブルはあらかじめ利用者がデータ加工の目的を考えあわせて、どのデータベースのどの属性を用いるかを設計して定義するものとする。

【0039】この後のルールインダクション処理のフローチャートを図15に示す。まず、ステップ340において、クライアント処理装置(図1参照)は、図8に示された形式で事例テーブルの一部を表示装置に表示する。利用者は表示装置上で、結論属性と説明属性を指定する(ステップ341)。ここで、一般的に考えて明らかに結論属性と関係がない属性は省いておく。この例では、「氏名」、「氏名コード」、「住所」などの属性は結論属性「クレジット査定」と関係がないと考えられるから予め省いておく。ここではそれ以外の全ての属性を説明属性として指定した場合を考える。

【0040】次のステップ342において、クライアント処理装置は、説明属性および結論属性に対して、ラベル付与処理に必要な情報である「ラベルコード対応表」を作成する。図16にこのラベルコード対応表の例を示す。図16(a)は家族情報データベースに対するラベルコード対応表360、図16(b)は金融関連情報データベースに対するラベルコード対応表370である。対応表の意味するところは図11で説明したものと大体同じであるが、キー属性を指定する項目が含まれていることが異なっている。それ以外では、ここまでの処理は第1実施例で説明したものと同様である。ここで、ルールインダクションに用いる属性がそれぞれどのデータベースの属性であるかは、図14の属性関係づけテーブルを参照して決定することができる。

【0041】次に利用者によってルールインダクション処理の開始が指示される(ステップ343)と、クライアント処理装置は、2つのサーバ処理装置に対し、それぞれ図16(a)、(b)に示したラベルコード対応表360、370からなる検索要求を送信する。

【0042】ステップ345において、各サーバ処理装置(図1参照)は図16の検索要求を受信すると、それぞれのサーバ処理装置において検索処理を開始する。まずそれぞれの事例データを、図16(a)、(b)のラベルコード対応表に基づいて、各事例を図6と同様の事例一属性値対応表に変換する(ステップ346)。図17にその一例を示す。この検索処理では、図16(a)および(b)のラベルコード対応表360および370に現れる全ての属性に対して、その属性が記号属性であ

れば属性値を指定されたラベルコードに置き換え、数値属性であればラベルコード対応表に指定された数値範囲で分類してラベルコードに置き換える処理を行う。ここで、第1の実施例における図12との違いは、複数のデータベースには事例テーブルが列方向に分割されて格納されているため、図17の事例－属性値対応表にはレコードに固有な識別情報として指定されたキー属性である「氏名コード」の情報が付加されていることである。各サーバ処理装置は、ローカルに管理している事例テーブルを図17の事例－属性値対応表に変換した後、ステップ347において、その結果を検索結果としてクライアント処理装置に送信する。

【0043】クライアント処理装置は、複数のサーバ処理装置から処理結果を受信し（ステップ348）、レコード識別情報である「氏名コード」情報をキーとして合成して、事例テーブル全体の事例－属性値対応表をクライアント処理装置上に作成する（ステップ349）。

【0044】ここで重要なのは、2つのデータベースの整合性は管理されていない場合があるので、どちらか一方の事例テーブルにしか存在しないレコードがありうることである。このような場合、合成された事例－属性値対応表の中で対応しない部分は欠損値を表す特別な値を埋めこむことにすれば、キー属性を基に事例テーブル全体に対する事例－属性値対応表を合成することができる。事例－属性値対応表が作成された後は、ルールインダクション処理はこの表だけを参照してルールを抽出する（ステップ350）ことは、第1実施例と全く同様である。

【0045】以上の説明は、ルールインダクションに限らず、ID3のような属性に基づいて分類の決定木を自動生成する処理にも全く同様に適用することができる。以上述べた本発明の第2実施例によれば、クライアント処理装置におけるルールインダクション処理の実行に当たって、あらかじめ複数のサーバ処理装置で実行されるDBMS（データベース管理システム）の機能を拡張しておき、クライアント処理装置で必要な検索結果をクライアント処理装置からの要求に従って並列に実行させることが可能になる。複数のサーバ処理装置への検索要求とその結果の取得は、ルールインダクション処理本体の実行に先立って一度だけ実行すればよく、ネットワーク上で流れるデータ量を大幅に削減し、ネットワーク負荷を減少させることが可能になる。また、検索処理自身も並列に実行されるため、全体の処理時間を短縮できるという効果もある。さらに利用者から見れば、あくまでクライアント処理装置上の表イメージのデータに対して直接処理を実行し、結果を得るという作業が可能になるため、事例テーブルが複数のデータベースに分割して格納されていることは全く意識することがなく、より自然に作業を進めることができるという効果もある。

【0046】以上の効果は第1実施例と同様であるが、

さらに第2実施例特有の効果としては、整合性の保証されていない複数の独立したデータベース上に分割して格納されている情報を、あたかも1つの仮想的なデータベースとして取り扱うことができるという効果もある。データ加工処理をこのように実現することにより、従来の統計解析処理や、学習によるニューロモデルまたはファジモデルの構築などと全く同様のインタフェースを用いてルールインダクションなどの処理を行うことができる。このように統一されたインタフェースでさまざまな種類のデータ加工処理を自由に実行し、結果を分析することができるので、利用者の作業効率が向上するという効果もある。

【0047】

【発明の効果】以上、本発明によれば複数のデータベースに分割して格納されているデータに対してルールインダクションなどのデータ加工処理を適用する際に、複数のデータベースが備えているデータ処理能力を並列に活用することができるので、高速な処理が可能になるという効果がある。また、データ加工に用いるデータそのものを所定の装置（クライアント処理装置）上に集める必要もなく、また複数のデータベースに対する検索処理はそれぞれ1回で十分なため、ネットワークの負荷を大幅に軽減することができるという効果がある。さらに、利用者にとっては全てテーブル形式で仮想的に定義されたデータを対象に処理を進めることができ、実際のデータが物理的に分割して複数のデータベースに格納されていることを意識する必要がないので、データ加工処理を容易に進めることができるという効果もある。

【図面の簡単な説明】

【図1】本発明の対象とする複数データベース利用データ加工処理システムの全体構成図である。

【図2】複数のデータベースへのデータ分割格納方式を説明するための図である。

【図3】事例テーブルの例を示す図である。

【図4】ラベル付与処理により記号属性に変換された事例テーブルの例を示す図である。

【図5】ルールインダクションにより抽出されたルールの例を示す図である。

【図6】事例－属性値対応表の例を示す図である。

【図7】表成分値－属性値対応表を示す図である。

【図8】クレジット査定における事例テーブルの構成図である。

【図9】事例テーブルの分割格納方式を示す図である。

【図10】本発明の第1実施例の処理フローチャートである。

【図11】ラベルコード対応表の例を示す図である。

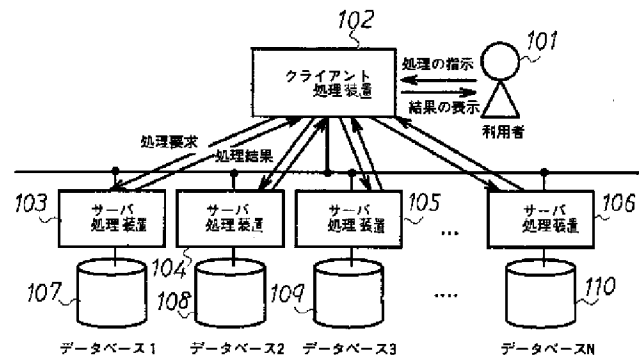
【図12】サーバ処理装置における事例－属性値対応表の例を示す図である。

【図13】列分割による事例テーブルの分割格納方式を示す図である。

【図14】属性関連づけテーブルの例を示す図である。
 【図15】本発明の第2実施例の処理フローチャートである。
 【図16】ラベルコード対応表の例を示す図である。
 【図17】サーバ処理装置における事例－属性値対応表の例を示す図である。
 【図18】ID3によって生成された分類木の例を示す

図である。
 【符号の説明】
 101 利用者
 102 クライアント処理装置
 103～106 サーバ処理装置
 107～110 データベース

【図1】



【図3】

201	202	203	204
銀行の信用	預金残高	現在の負債	クレジット査定
有	400	2200	可
有	1000	2500	可
無	2000	1000	可
有	-2000	10000	不可
有	400	6000	可
有	100	1000	可
無	400	8000	不可
無	1000	5000	不可
無		9500	不可

【図4】

【図2】

	属性1	属性2	属性3	...	属性K	
1						
2						
3						
4						
...						
M						

レコード番号

→ データベース1に格納
 → データベース2に格納
 → データベース3に格納

207	208	209	210
銀行の信用	預金残高	現在の負債	クレジット査定
有	中	小	可
有	大	小	可
無	大	小	可
有	小	大	不可
有	中	中	可
有	中	小	可
無	中	大	不可
無	大	中	不可
無	—	大	不可

【図5】

ルール1
 IF 現在の負債=小
 THEN 預金残高=中または大
 THEN クレジット査定=可

ルール2
 IF 現在の負債=大
 THEN クレジット査定=不可

【図6】

212	213	214	215	216
レコード番号	銀行の信用	預金残高	現在の負債	クレジット査定
1	1	1	0	1
2	1	2	0	1
3	0	2	0	1
4	1	0	2	0
5	1	1	1	1
6	1	1	0	1
7	0	1	2	0
8	0	2	1	0
9	0	3	2	0

【図 7】

(a)

無	0
有	1

(b)

小	0
中	1
大	2
不明	3

(c)

小	0
中	1
大	2

(d)

不可	0
可	1

【図 8】

221 222 223 224 225 226 227 220

事例番号	氏名	氏名コード	信用	残高	年次	負債
1	青柳一郎	0100235	有	1200	900	28.00
2	鈴木宏	7900015	無	200	1200	220
3	加藤氏一	7601002	無	550	1350	50
4	飯田裕子	6600105	有	100	250	90
5	清水健一	9000115	有	780	1370	
6	広沢孝司	7701059	無	2100	1100	680

228 229 230 231 232 233 234 235 236 237 238

年齢	性別	住所	配偶者	扶養家族	職業	最終学歴	住居	カード枚数	利用実績	査定
39	男	東京都	有	3	会社員	大学	自宅	0	無	可
41	男	大阪市	有	2	会社員	大学	社宅	2	小	可
41	男	福岡市	有	2	自営業	高校	一	3	中	不可
34	女	岡山市	無	0	無職	短大	自宅	0	無	不可
30	男	横浜市	無	1	学生	高校	自宅	1	中	可
43	男	横浜市	有	4	会社員	大学院	借家	5	大	可

【図 9】

【図 10】

240

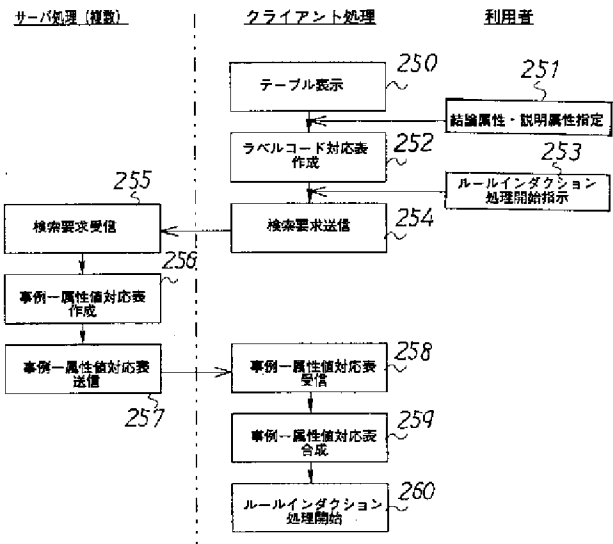
年月	氏名	コード	信用	扶養	年次	負債	不動産	性別	住所	配偶者
1										
2										
...										
1000										
1001										
1002										
...										
2500										
2501										
2502										
...										
4500										
4501										
...										

扶養家族	職業	学歴	住居	カード枚数	利用実績	査定

→データベース 1

→データベース 2

→データベース 3



【図 11】

261

属性値	ラベル値	最小値	最大値
信用：無	0	-	-
信用：有	1	-	-
残高：小	0	-	0
残高：中	1	0	1000
残高：大	2	1000	-
残高：不明	3	-	-
年齢：若い	0	0	20
年齢：青年	1	20	35
年齢：中年	2	35	50
年齢：壮年	3	50	-
性別：男	0	-	-
性別：女	1	-	-
:	:	:	:

264

262

263

265

【図12】

レコード番号	信用	残高	年齢	性別	配偶者	...
1	1	2	1	0	1	
2	0	1	2	0	1	
3	0	1	2	0	1	
4	1	1	0	?	0	
5	1	3	1	1	1	
6	0	2	1	0	1	
...

【図14】

属性	データベース種別	項目名
氏名	家族情報DB	氏名
氏名コード	キー属性	氏名コード
銀行の信用	金融関係情報DB	信用
残高	金融関係情報DB	預金残高
年収	家族情報DB	年収
...

【図16】

属性値	ラベル値	最小値	最大値
キー属性：氏名コード			
年齢：若い	0	0	20
年齢：青年	1	20	35
年齢：中年	2	35	50
年齢：壮年	3	50	-
性別：男	0	-	-
性別：女	1	-	-
...

属性値	ラベル値	最小値	最大値
キー属性：氏名コード			
残高：小	0	-	0
残高：中	1	0	1000
残高：大	2	1000	-
残高：不明	3	-	-
...

【図17】

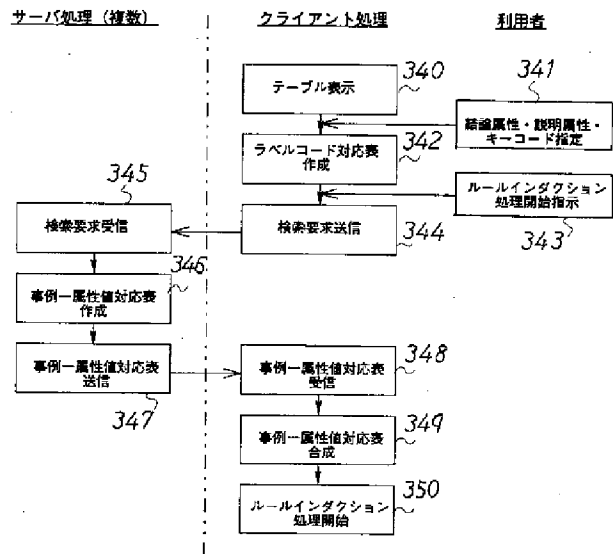
氏名コード	レコード番号	信用	残高	年齢	性別	配偶者	...
8100235	1	1	2	1	0	1	
7600215	2	0	1	2	0	1	
7601002	3	0	1	2	0	1	
8600306	4	1	1	0	?	0	
8600116	5	1	3	1	1	1	
7701059	6	0	2	1	0	1	
...

【図13】

氏名	コード	信用	残高	年収	職業	カード枚数	利用実績	審査
鈴木一郎	8100235	無	1200	950	会社員	0	無	可
...

氏名	コード	年齢	性別	住所	配偶者	扶養家族	職業	年収	学歴	住居
鈴木一郎	8100235	39	男	東京都	有	3	会社員	900	大学	自主
...

【図15】



【図18】

